

 Open access • Posted Content • DOI:10.1101/384016

Standardized biogeographic grouping system for annotating populations in pharmacogenetic research — [Source link](#)

Rachel Huddart, [Alison E. Fohner](#), [Michelle Whirl-Carrillo](#), [Genevieve L. Wojcik](#) ...+5 more authors

Institutions: [Stanford University](#), [University of Washington](#), [Anschutz Medical Campus](#)

Published on: 03 Aug 2018 - [bioRxiv](#) (Cold Spring Harbor Laboratory)

Topics: [Genetic variability](#), [Population](#), [1000 Genomes Project](#) and [Genetic structure](#)

Related papers:

- [Genetic Similarities Within and Between Human Populations](#)
- [Detecting genetic isolation in human populations: a study of European language minorities.](#)
- [A variant-centric perspective on geographic patterns of human allele frequency variation](#)
- [Demographic history and rare allele sharing among human populations](#)
- [Limited Evidence for Classic Selective Sweeps in African Populations](#)

Share this paper:    

View more about this paper here: <https://typeset.io/papers/standardized-biogeographic-grouping-system-for-annotating-5030p9smrj>

1 **Title: Standardized biogeographic grouping system for annotating populations in**
2 **pharmacogenetic research**

3
4 **Authors:**

5 Rachel Huddart^{1,8}, Alison E. Fohner^{1,2,8}, Michelle Whirl-Carrillo^{1,8}, Genevieve L. Wojcik¹,
6 Christopher R. Gignoux³, Alice B. Popejoy^{1,4}, Carlos D. Bustamante^{1,5}, Russ B. Altman^{1,5,6,7},
7 Teri E. Klein^{1,7*}

8
9 **Affiliations:**

10 ¹ Department of Biomedical Data Science, Stanford University, Stanford, CA 94305, USA

11 ² Department of Epidemiology, University of Washington, Seattle, WA 98195, USA

12 ³ Division of Bioinformatics and Personalized Medicine, and Department of Biostatistics,
13 University of Colorado, Anschutz Medical Campus, Aurora, CO 80045, USA

14 ⁴ Stanford Center for Integration of Research on Genetics and Ethics, Stanford, CA 94305, USA

15 ⁵ Department of Genetics, Stanford University, Stanford, CA 94305, USA

16 ⁶ Department of Biomedical Engineering, Stanford University, Stanford, CA 94305, USA

17 ⁷ Department of Medicine, Stanford University, Stanford, CA 94305, USA

18 ⁸ These authors contributed equally to this work

19 * **Correspondence:** Dr. Teri E. Klein, Shriram Center for BioE & ChemE, 443 Via Ortega,

20 Stanford, CA 94305-4125, Phone: (650) 736-0156, feedback@pharmgkb.org

21 **Abstract:**

22 The varying frequencies of pharmacogenetic alleles between populations have important
23 implications for the impact of these alleles in different populations. Current population grouping
24 methods to communicate these patterns are insufficient as they are inconsistent and fail to reflect
25 the global distribution of genetic variability. To facilitate and standardize the reporting of
26 variability in pharmacogenetic allele frequencies, we present seven geographically-defined
27 groups: American, Central/South Asian, East Asian, European, Near Eastern, Oceanian, and
28 Sub-Saharan African, and two admixed groups: African American/Afro-Caribbean and Latino.
29 These nine groups are defined by global autosomal genetic structure and based on data from
30 large-scale sequencing initiatives. We recognize that broadly grouping global populations is an
31 oversimplification of human diversity and does not capture complex social and cultural identity.
32 However, these groups meet a key need in pharmacogenetics research by enabling consistent
33 communication of the scale of variability in global allele frequencies and are now used by
34 PharmGKB.

35

36 **Introduction**

37 Interindividual variability in pharmacogenes has important consequences for drug
38 efficacy and toxicity.(1, 2) Unlike the low frequencies of alleles that are considered actionable
39 with respect to disease risk, pharmacogenetic variants with clinical relevance are common and,
40 in fact, both presence and absence of variants provide valuable dosing information.(3, 4) The
41 frequencies of many pharmacogenetic alleles vary greatly by global population, meaning that
42 people with different ancestries can have considerably different likelihoods of carrying an allele
43 that is associated with a particular drug response. For example, the CYP3A5*3 allele has been
44 found at a frequency of 98% in an Iranian population but at 11% in a Ngoni population from
45 Malawi. (5, 6) A single value for global allele frequency would fail to reflect this pattern.
46 Presenting the differences in frequencies of pharmacogenetic alleles is important for
47 communicating the scale of their expected impact on drug response and the degree of variation
48 between populations. This information is invaluable for furthering pharmacogenetic research and
49 implementation.

50 Many pharmacogenetic studies present allelic data for very specific populations, such as
51 from a single country or ethnic group, which are difficult to incorporate into broader research or
52 implementation. Literature curation and gene summaries, such as those from the
53 Pharmacogenomics Knowledgebase (PharmGKB: www.pharmgkb.org), must group these
54 specific populations when annotating pharmacogenetic studies to allow users to easily compare
55 information from multiple studies. As such, tagging studies with population group identifiers is
56 an important component of knowledge extraction from curated literature. These population group
57 labels then are used in aggregating and evaluating overall evidence for gene-drug associations,

58 which eventually inform clinical implementation guidelines, such as those of the Clinical
59 Pharmacogenetics Implementation Consortium (CPIC: www.cpicpgx.org).

60 Similar to other areas of biomedical research, (7) current methods for grouping global
61 populations in pharmacogenetics are based on subjective, vague, and inconsistent geographical
62 boundaries, or on populations that are geographically straightforward to cluster and reflect little
63 admixture.(8-12) As an example of the issues with current grouping methods, some studies
64 cluster participants of Egyptian descent with African populations, while others cluster them with
65 Middle Eastern populations.(13, 14) While this discrepancy illustrates inconsistencies of
66 geographic borders, the clustering of African-descent populations of the Americas with
67 populations from Africa, as seen in the 1000 Genomes African (AFR) superpopulation, provides
68 another example of challenges posed by employing a small number of categories to describe a
69 broad spectrum of genomically diverse groups. The genetic patterns seen in American
70 populations with African ancestry differs dramatically from populations in Africa due to
71 admixture primarily with European and American Indian populations. (15-17) While sharing
72 common ancestry, the recent admixture typically observed in the Americas can complicate
73 average allele frequency estimation or, at a minimum, make these combined groupings less
74 homogeneous.(16) These insufficient grouping systems, often ad-hoc and not fully representative
75 evidence from population genomic studies, create a barrier to understanding and interpreting
76 pharmacogenetic allele frequencies in a globally representative fashion.

77 Until July 2018, PharmGKB annotated studies using the five race categories defined by
78 the US Office of Management and Budget (OMB): White, Black or African American, American
79 Indian or Alaska Native, Asian, and Native Hawaiian or Pacific Islander, with an additional
80 ethnicity OMB category of Hispanic/Latino. While PharmGKB serves as a global resource, these

81 OMB groups are US-centric and, as socio-cultural measures of identity, lack the capacity to
82 capture the scale of global human diversity. We also investigated the utility of the biogeographic
83 categories employed by the Human Genome Diversity Panel - Centre d'Etude du Polymorphisme
84 Humain (HGDP - CEPH), which groups its 52 populations into Africa, Europe, Middle East,
85 South and Central Asia, East Asia, Oceania and the Americas.(8, 18, 19) These population labels
86 work well for the populations included in the HGDP data set, which are not located in
87 ambiguous regions between group borders and which mostly contain populations with little
88 admixture. However, papers curated at PharmGKB can include populations located all over the
89 world, including in the transitional zones between HGDP geographical regions and admixed
90 populations. This leads to ambiguity in how such populations would be grouped using HGDP
91 categories. In conclusion, existing systems are insufficient for capturing the diversity of study
92 populations in a replicable manner that is consistent with patterns of human genetic variation.

93 Therefore, we sought to define a grouping system of global populations that could be
94 used consistently to annotate pharmacogenetic studies and relevant alleles, and could capture
95 global human population genetic patterns. Using population genetics data sources, including the
96 1000 Genomes Phase 3 data release and the HGDP, we propose a simple and robust grouping
97 pattern based on nine broad biogeographic regions that represent major geographic regions of the
98 world (**Figure 1**). It is important to note that classifying individuals and communities into a few
99 distinct groups with defined boundaries conflicts with our understanding of human variation,
100 history, and social/cultural identities. *As a result, we respectfully present these groups as a tool*
101 *to represent broad differences in frequencies of pharmacogenetic variation rather than as a*
102 *classification of human diversity.*

103

104

105 **Results**

106 We chose this geographic clustering pattern because geography has historically been the
107 greatest predictor of genetic variation between human populations, with genetic distance
108 increasing as geographic distance increases.(20) This geographic pattern aids consistency in
109 population groupings by setting boundaries along national borders. To simplify utility,
110 geographic boundaries between groupings are drawn predominantly along country borders, with
111 only Russia divided into east and west along the Ural Mountains boundary due to the large size
112 and genetic heterogeneity of the country. We intend these groups to represent peoples with a
113 predominance of ancestors who were in the region pre-Diaspora and pre-colonization.

114 We have also included two admixed groups representing populations with recent gene
115 flow between geographically-based populations and therefore, have distinct genetic patterns
116 which are not adequately reflected by any single geographically-based group. (7) While many
117 populations reflect a degree of admixture, we selected these two populations because they are
118 frequently reported in pharmacogenetic studies.

119 We consider these nine groups sufficient to better illustrate the broad diversity in global
120 allele frequencies, yet small enough to apply easily and to be tractable in grouping specific
121 populations.(21-24) The groups are given below with their abbreviations.

122

123 **Geographical populations**

124 *American (AME)*: The American genetic ancestry group includes populations from both North
125 and South America with ancestors predating European colonization, including American Indian,

126 Alaska Native, First Nations, Inuit, and Métis in Canada, and Indigenous peoples of Central and
127 South America.

128 *Central/South Asian (SAS)*: The Central and South Asian genetic ancestry group includes
129 populations from Pakistan, Sri Lanka, Bangladesh, India, and ranges from Afghanistan to the
130 western border of China.

131 *East Asian (EAS)*: The East Asian genetic ancestry group includes populations from Japan,
132 Korea, and China, and stretches from mainland Southeast Asia through the islands of Southeast
133 Asia. In addition, it includes portions of central Asia and Russia east of the Ural Mountains.

134 *European (EUR)*: The European genetic ancestry group includes populations of primarily
135 European descent, including European Americans. We define the European region as extending
136 west from the Ural Mountains and south to the Turkish and Bulgarian border.

137 *Near Eastern (NEA)*: The Near Eastern genetic ancestry group encompasses populations from
138 northern Africa, the Middle East, and the Caucasus. It includes Turkey and African nations north
139 of the Saharan Desert.

140 *Oceanian (OCE)*: The Oceanian genetic ancestry group includes pre-colonial populations of the
141 Pacific Islands, including Hawaii, Australia, and Papua New Guinea.

142 *Sub-Saharan African (SSA)*: The Sub-Saharan African genetic ancestry group includes
143 individuals from all regions in Sub-Saharan Africa, including Madagascar.(25)

144

145 **Admixed populations**

146 *African American/Afro-Caribbean (AAC)*: Individuals in the African American/Afro-Caribbean
147 genetic ancestry group reflect the extensive admixture between African, European, and
148 Indigenous ancestries(26) and, as such, display a unique genetic profile compared to individuals

149 from each of those lineages alone. Examples within this cluster include the Coriell Institute's
150 African Caribbean in Barbados (ACB) population and the African Americans from the
151 Southwest US (ASW) population, (27) and individuals from Jamaica and the US Virgin Islands.
152 *Latino (LAT)*: The Latino genetic ancestry group is not defined by an exclusive geographic
153 region, but includes individuals of Mestizo descent, individuals from Latin America, and self-
154 identified Latino individuals in the United States. Like the African American/Afro-Caribbean
155 group, the admixture in this population creates a unique genetic pattern compared to any of the
156 discrete geographic regions, with individuals reflecting mixed Native and Indigenous American,
157 European, and African ancestry.

158

159 The Central/South Asian, East Asian and European groups presented here are equivalent
160 to the 1000 Genomes South Asian (SAS), East Asian (EAS) and European (EUR) super
161 populations, respectively. As such, we have adopted the relevant 1000 Genomes super
162 population codes as abbreviations for each of these groups to maintain consistency. While the
163 1000 Genomes Ad Mixed American (AMR) super population shows complete overlap with the
164 Latino group, we have opted to use the abbreviation LAT for this group. This removes the
165 potential for confusion between the Latino group and the other admixed group of African
166 American/Afro-Caribbean.

167 **Figure 1** illustrates the countries included in each of the seven geographical groups and
168 removes any ambiguity of the group boundaries. As this map shows the boundaries of each
169 group pre-colonization and pre-Diaspora, the two admixed groups, African American/Afro-
170 Caribbean and Latino are not shown. We intend this map to be used as a guide for grouping
171 genetic ancestral populations. Study subjects of an ancestry that is not within the geographic

172 cluster in which they currently live will be included in the geographic cluster reflecting their
173 ancestry. For example, South Africans of Dutch descent would be included in the European
174 cluster rather than the Sub-Saharan African cluster. However, when lacking a clear description
175 otherwise, the population will be included in the group that includes its home country.

176 This approach highlights the importance of understanding and recording detailed self-
177 identified and self-reported race and ethnicity in the context of genetic studies. While self-
178 reported race and ethnicity can be influenced by an individual's social and cultural background
179 and thus may not perfectly correlate with genetic ancestry (28), it is more reliable than
180 assignment of race or ethnicity by another person (e.g. a healthcare professional) (29). However,
181 it should be noted that self-reported measures can be complicated by collection processes, (30)
182 including an incomplete selection of possible identity categories, or allowing only one selection
183 and thus failing to capture whether an individual may identify with multiple categories or none at
184 all (29). These classification limitations can be particularly prevalent among populations with a
185 high degree of admixture.

186 To validate the genetic variability distinguished by these population groups, we
187 conducted Principal Components Analysis (PCA) using autosomal genotype data of unrelated
188 individuals from 1000 Genomes and HGDP. As seen in **Figure 2A**, the first two principal
189 components (PCs) separate populations by geographic region, especially along continental
190 boundaries, and illustrate the increasing genetic distance between populations of increasing
191 geographic distance. As can be seen in the overlapping PC distribution of individuals of different
192 population groups, human genetic diversity is a spectrum,(19) and therefore the geographic
193 boundaries of these groups should be understood as an obligatory divide to create relevant
194 groupings, with the acknowledgement that these borders are constrained by modern country

195 borders and therefore are inherently arbitrary in geographic space.(19) However, as shown in
196 **Figure 2B**, only a few PCs are needed to accurately predict these population clusters. Even with
197 only 4 PCs, the minimum area under the curve (AUC) for correct cluster prediction is 97.9% for
198 most populations using multiple logistic regression. The only outlier is the African
199 American/Afro-Caribbean cluster, consistent with ancestral similarity to the African cluster.(15,
200 31) Here still, with a larger number of PCs, the AUC is above 93%, even with the observed
201 ancestry outliers present in the 1000 Genomes African Americans in the Southwest US (ASW)
202 population.(32) While no categorization will result in perfect prediction, given the spectrum of
203 human diversity, the statistical validation of this clustering from broad autosomal data makes
204 these clusters both relevant and useful for PharmGKB.

205
206 In **Figure 3**, we demonstrate that the groups we have selected are effective for
207 representing the diversity of global allele frequencies in pharmacogenes. We present here the
208 frequency of four single nucleotide polymorphisms (SNPs) with important pharmacogenetic
209 implications. The ‘A’ allele of rs1065852 is the defining SNP of the *cytochrome P450 2D6*
210 (*CYP2D6*) *10 haplotype and is also found in combination with other variants in multiple
211 CYP2D6 haplotypes. Haplotypes containing this SNP are associated with decreased CYP2D6
212 activity, which has important implications for drugs that are CYP2D6 substrates, including
213 codeine, selective serotonin reuptake inhibitors, ondansetron, and tricyclic antidepressants.(33-
214 36) The *CYP2C9* alleles *2 (defined by rs1799853), *3 (defined by rs1057910), and *8 (defined
215 by rs7900194) are associated with reduced enzyme function and therefore are associated with
216 recommended changes to the dosing of warfarin and phenytoin, which are substrates of
217 CYP2C9.(37, 38) Using data from the 1000 Genomes, we show the frequency of the four SNPs

218 in these biogeographic groups. The range of frequencies between populations illustrates the
219 importance of showing allele frequency by group in order to convey its impact on drug response
220 globally.

221 The SNP rs1065852 shows stark continental patterns (**Figure 3A**). The ‘A’ allele is found
222 at high frequencies within East Asian populations, ranging from 66.2% in Vietnam (KHV) to
223 36.1% in Japan (JPT). This allele is less frequent in other continental populations, such as Sub-
224 Saharan African (3.5-16.5%), European (14.6-24.7%), and Central/South Asian (10.4-25.6%).
225 As can be seen from the range of frequencies of the three *CYP2C9* alleles, the most common
226 reduced function allele varies globally, with the *8 allele much more common in Sub-Saharan
227 African populations (1.8-7.6%) than the *2 (<1%) or *3 (monomorphic in Africa) (**Figure 3B-**
228 **D**). Conversely, the *8 allele is rare in European populations (<1%), while *2 (8.1-15.2%) and
229 *3 (5.6-8.4%) are more common. Patterns such as this one can result in bias in the utility of
230 dosing algorithms, such as the International Warfarin Pharmacogenetics Consortium (IWPC)
231 dosing algorithm for warfarin, which adjusts dose based on the presence of the *2 and *3 alleles
232 but does not include the *8 allele.(39)

233

234 **Discussion**

235 While individual pharmacogenetic testing (either pre-emptive or at point-of-care) remains
236 the most effective and appropriate way to implement pharmacogenetic knowledge for the care of
237 an individual,(40, 41) we recognize the need in clinical and genetic research for a standardized
238 method to broadly group populations based on biogeographic region. For example, identifying
239 populations with high frequencies of certain pharmacogenetic alleles can help to direct targeted
240 screening when resources are constrained and inform priorities for future pharmacogenetic

241 research.(20) However, the groups we present are large and the summary information presented
242 should be understood as an approximation dependent on existing studies in that region, which
243 may be limited to a few locations. As such, these groups are not suitable for use in guiding
244 specific implementation programs; rather, they should be seen as a tool for research purposes.

245 It should be noted that this grouping system does have limitations. Classifying
246 individuals into these population groups can be complicated by social and cultural identities(8,
247 10, 42-44) and membership of an individual within one of these population groups is inherently
248 an imperfect surrogate for predicting the likelihood that the individual carries a particular genetic
249 variant.(41, 45) As can be seen in the analysis of rs1065852 above, the frequency of the ‘A’
250 allele can vary by up to 30% between populations which are all included in the East Asian group.
251 Furthermore, while the grouping system is based on overall genome-wide average patterns,
252 which typically follow a clinal variation pattern correlated with geographic proximity,(8, 23, 24,
253 46, 47) variation in individual genes or individual populations do not always follow these
254 gradual patterns.(9-12, 41) In an attempt to mitigate some of these limitations, we encourage
255 researchers using this grouping system to also provide specific details regarding the geographical
256 and racial or ethnic origins of their subjects.

257 Because aggregate annotations of pharmacogenetic research and summary allele
258 frequencies are based only on available studies, additional studies are needed that include a
259 greater diversity of populations to make pharmacogenetic research and allele frequency
260 summaries more representative.(48) For example, the Sub-Saharan African (SSA) grouping
261 represents a large swath of human genomic diversity, which is not adequately represented in the
262 available data from HGDP and 1000 Genomes. Increased representation of these populations in
263 pharmacogenetics studies may lead to the discovery of clinical differences within the larger

264 grouping. Furthermore, large, reference genetic studies with targeted allele information, like that
265 emerging from the Population Architecture using Genomics and Epidemiology (PAGE) study
266 (www.pagestudy.org), may provide compelling evidence to adjust these group boundaries based
267 on frequency patterns specific to pharmacogenetic alleles. Continued evolution of this grouping
268 system will be key to ensuring that misclassification of individuals is kept to a minimum.
269 However, it should be understood that some misclassification is inevitable and will only be truly
270 avoided when every patient can access comprehensive pharmacogenetic testing.

271 Despite these limitations, broad population groups are needed for illustrating global
272 diversity with respect to pharmacogenetic variation and the average predicted phenotypes in
273 populations. These nine proposed biogeographic groups provide a consistent way to present
274 these data based on a system that is grounded in robust data on population genetic patterns, and
275 their introduction is particularly timely given the recent commentaries by Bonham *et al.* and
276 Cooper *et al.* (7, 49) PharmGKB is now using these population groups in curation activities, and
277 we recommend that these groups and accompanying map be considered the standard grouping
278 mechanism for population pharmacogenetics. Ultimately, individual pharmacogenetic testing of
279 all patients, regardless of ancestry, is needed to deliver truly personalized medicine. However,
280 the population groups we present are useful for the standardized presentation of pharmacogenetic
281 studies, global allele frequency summaries in pharmacogenetic research and broad clinical
282 screening.

283

284 **Methods**

285 The MVN joint callset for 1000 Genomes data Phase 3 (21) was downloaded directly from the
286 website for downstream interpretation. For principal component analysis (PCA), we filtered sites

287 with a MAF < 0.5% and thinned sites given windows of 100 kilobases or 10 variants and $r^2 > 0.2$,
288 resulting in 156,211 sites. PCA was performed in PLINK 1.9 (50). Forward stepwise logistic
289 regression was subsequently performed, adding 1 PC at a time, to predict population labels in a
290 bivariate fashion. Prediction accuracy was assessed using the AUC-ROC estimator, as included
291 in the R package ‘epicalc.’ To make assessments transparent, we included all individuals with
292 specific population labels, although it has been demonstrated in multiple venues that there are
293 several known ancestry outliers within 1000 Genomes populations of the Americas (17, 32).
294 Plots were performed in R and ggplot2.

295

296 **Study Highlights**

297 *What is the current knowledge on the topic?*

298 The frequency of pharmacogenetic alleles can vary significantly between different populations
299 around the world. Grouping populations can simplify reporting of pharmacogenetic alleles but
300 current methods used to group populations are inadequate and are applied inconsistently.

301 *What question did this study address?*

302 Can we improve how populations are grouped for the reporting of pharmacogenetic alleles?

303 *What does this study add to our knowledge?*

304 We present nine new biogeographical groups based on geographical location or recent genetic
305 admixture for use in pharmacogenetic research. These groups have been validated using
306 autosomal genetic data from large-scale sequencing initiatives.

307 *How might this change clinical pharmacology or translational science?*

308 These groups have already been adopted for use in curation activities at PharmGKB. It is hoped
309 that use of these groups will become standard in pharmacogenetics research.

310

311 **Author Contributions**

312 R.H., A.E.F., M.W-C., G.L.W., C.R.G., A.B.P., C.D.B., R.B.A. and T.E.K. wrote the manuscript;

313 A.E.F., M.W-C., C.R.G. and T.E.K. designed the research, M.W-C., G.L.W. and C.R.G. analyzed

314 the data.

315 **Conflicts of Interest:**

316 CRG owns stock in 23andMe, Inc and is a founder of and advisor to Encompass Bioscience, Inc.

317 CDB is a member of the scientific advisory boards for Liberty Biosecurity, Personalis, 23andMe

318 Roots into the Future, Ancestry.com, IdentifyGenomics, and Etalon and is a founder of CDB

319 Consulting. RBA is a stockholder in Personalis Inc. and 23andMe, and a paid advisor for

320 Youscript. Remaining authors have no conflicts of interest.

321

322 **References:**

- 323 1. Roden, D.M. (2006). PHarmacogenomics: Challenges and Opportunities. *Annals of Internal*
324 *Medicine* 145.
- 325 2. Dunnenberger, H.M., Crews, K.R., Hoffman, J.M., Caudle, K.E., Broeckel, U., Howard, S.C.,
326 Hunkler, R.J., Klein, T.E., Evans, W.E., and Relling, M.V. (2015). Preemptive clinical
327 pharmacogenetics implementation: current programs in five US medical centers. *Annu*
328 *Rev Pharmacol Toxicol* 55, 89-106.
- 329 3. Tabor, H.K., Auer, P.L., Jamal, S.M., Chong, J.X., Yu, J.H., Gordon, A.S., Graubert, T.A.,
330 O'Donnell, C.J., Rich, S.S., Nickerson, D.A., et al. (2014). Pathogenic variants for
331 Mendelian and complex traits in exomes of 6,517 European and African Americans:
332 implications for the return of incidental results. *Am J Hum Genet* 95, 183-193.
- 333 4. Wright, G.E.B., Carleton, B., Hayden, M.R., and Ross, C.J.D. (2018). The global spectrum of
334 protein-coding pharmacogenomic diversity. *The pharmacogenomics journal* 18, 187-195.
- 335 5. Rahsaz, M., Azarpira, N., Nikeghbalian, S., Aghdaie, M.H., Geramizadeh, B., Moini, M.,
336 Banihashemi, M., Darai, M., Malekpour, Z., and Malekhosseini, S.A. (2012). Association
337 between tacrolimus concentration and genetic polymorphisms of CYP3A5 and ABCB1
338 during the early stage after liver transplant in an Iranian population. *Experimental and*
339 *clinical transplantation : official journal of the Middle East Society for Organ*
340 *Transplantation* 10, 24-29.
- 341 6. Bains, R.K., Kovacevic, M., Plaster, C.A., Tarekegn, A., Bekele, E., Bradman, N.N., and
342 Thomas, M.G. (2013). Molecular diversity and population structure at the Cytochrome
343 P450 3A5 gene in Africa. *BMC genetics* 14, 34.

- 344 7. Rosenberg, N.A., Pritchard, J.K., Weber, J.L., Cann, H.M., Kidd, K.K., Zhivotovsky, L.A., and
345 Feldman, M.W. (2002). Genetic structure of human populations. *Science (New York,*
346 *NY)* 298, 2381-2385.
- 347 8. Rajagopalan, R., and Fujimura, J.H. (2012). Will personalized medicine challenge or reify
348 categories of race and ethnicity? *The virtual mentor : VM* 14, 657-663.
- 349 9. Gannett, L. (2005). Group Categories in Pharmacogenetics Research. *Philosophy of Science*
350 72, 1232-1247.
- 351 10. Wilson, J.F., Weale, M.E., Smith, A.C., Gratrix, F., Fletcher, B., Thomas, M.G., Bradman, N.,
352 and Goldstein, D.B. (2001). Population genetic structure of variable drug response. *Nat*
353 *Genet* 29, 265-269.
- 354 11. Race, E., and Genetics Working Group. (2005). The Use of Racial, Ethnic, and Ancestral
355 Categories in Human Genetics Research. *Am J Hum Genet* 77, 519-532.
- 356 12. Relling, M.V., Gardner, E.E., Sandborn, W.J., Schmiegelow, K., Pui, C.H., Yee, S.W., Stein,
357 C.M., Carrillo, M., Evans, W.E., and Klein, T.E. (2011). Clinical Pharmacogenetics
358 Implementation Consortium guidelines for thiopurine methyltransferase genotype and
359 thiopurine dosing. *Clin Pharmacol Ther* 89, 387-391.
- 360 13. Scott, S.A., Sangkuhl, K., Stein, C.M., Hulot, J.S., Mega, J.L., Roden, D.M., Klein, T.E.,
361 Sabatine, M.S., Johnson, J.A., and Shuldiner, A.R. (2013). Clinical Pharmacogenetics
362 Implementation Consortium guidelines for CYP2C19 genotype and clopidogrel therapy:
363 2013 update. *Clin Pharmacol Ther* 94, 317-323.
- 364 14. Bryc, K., Auton, A., Nelson, M.R., Oksenberg, J.R., Hauser, S.L., Williams, S., Froment, A.,
365 Bodo, J.M., Wambebe, C., Tishkoff, S.A., et al. (2010). Genome-wide patterns of
366 population structure and admixture in West Africans and African Americans.
367 *Proceedings of the National Academy of Sciences of the United States of America* 107,
368 786-791.
- 369 15. Mathias, R.A., Taub, M.A., Gignoux, C.R., Fu, W., Musharoff, S., O'Connor, T.D., Vergara,
370 C., Torgerson, D.G., Pino-Yanes, M., Shringarpure, S.S., et al. (2016). A continuum of
371 admixture in the Western Hemisphere revealed by the African Diaspora genome. *Nature*
372 *communications* 7, 12522.
- 373 16. Martin, A.R., Gignoux, C.R., Walters, R.K., Wojcik, G.L., Neale, B.M., Gravel, S., Daly, M.J.,
374 Bustamante, C.D., and Kenny, E.E. (2017). Human Demographic History Impacts
375 Genetic Risk Prediction across Diverse Populations. *Am J Hum Genet* 100, 635-649.
- 376 17. Cann, H.M., de Toma, C., Cazes, L., Legrand, M.F., Morel, V., Piouffre, L., Bodmer, J.,
377 Bodmer, W.F., Bonne-Tamir, B., Cambon-Thomsen, A., et al. (2002). A human genome
378 diversity cell line panel. *Science (New York, NY)* 296, 261-262.
- 379 18. Rosenberg, N.A., Mahajan, S., Ramachandran, S., Zhao, C., Pritchard, J.K., and Feldman,
380 M.W. (2005). Clines, clusters, and the effect of study design on the inference of human
381 population structure. *PLoS Genet* 1, e70.
- 382 19. Burchard, E.G., Ziv, E., Coyle, N., Gomez, S.L., Tang, H., Karter, A.J., Mountain, J.L.,
383 Perez-Stable, E.J., Sheppard, D., and Risch, N. (2003). The importance of race and
384 ethnic background in biomedical research and clinical practice. *The New England journal*
385 *of medicine* 348, 1170-1175.
- 386 20. Auton, A., Brooks, L.D., Durbin, R.M., Garrison, E.P., Kang, H.M., Korbel, J.O., Marchini,
387 J.L., McCarthy, S., McVean, G.A., and Abecasis, G.R. (2015). A global reference for
388 human genetic variation. *Nature* 526, 68-74.
- 389 21. Elhaik, E., Tatarinova, T., Chebotarev, D., Piras, I.S., Maria Calo, C., De Montis, A., Atzori,
390 M., Marini, M., Tofanelli, S., Francalacci, P., et al. (2014). Geographic population
391 structure analysis of worldwide human populations infers their biogeographical origins.
392 *Nature communications* 5, 3513.

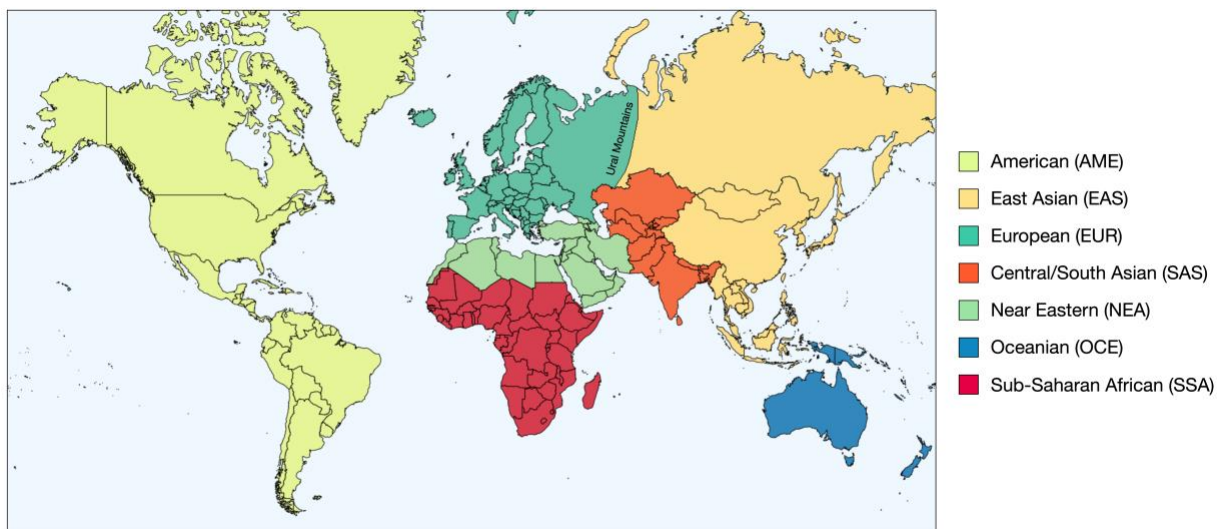
- 393 22. Jakobsson, M., Scholz, S.W., Scheet, P., Gibbs, J.R., VanLiere, J.M., Fung, H.C., Szpiech,
394 Z.A., Degnan, J.H., Wang, K., Guerreiro, R., et al. (2008). Genotype, haplotype and
395 copy-number variation in worldwide human populations. *Nature* 451, 998-1003.
- 396 23. Li, J.Z., Absher, D.M., Tang, H., Southwick, A.M., Casto, A.M., Ramachandran, S., Cann,
397 H.M., Barsh, G.S., Feldman, M., Cavalli-Sforza, L.L., et al. (2008). Worldwide human
398 relationships inferred from genome-wide patterns of variation. *Science* (New York, NY)
399 319, 1100-1104.
- 400 24. Hurles, M.E., Sykes, B.C., Jobling, M.A., and Forster, P. (2005). The dual origin of the
401 Malagasy in Island Southeast Asia and East Africa: evidence from maternal and paternal
402 lineages. *Am J Hum Genet* 76, 894-901.
- 403 25. Maples, B.K., Gravel, S., Kenny, E.E., and Bustamante, C.D. (2013). RFMix: a
404 discriminative modeling approach for rapid and robust local-ancestry inference. *Am J*
405 *Hum Genet* 93, 278-288.
- 406 26. Genomes Project, C., Abecasis, G.R., Auton, A., Brooks, L.D., DePristo, M.A., Durbin, R.M.,
407 Handsaker, R.E., Kang, H.M., Marth, G.T., and McVean, G.A. (2012). An integrated map
408 of genetic variation from 1,092 human genomes. *Nature* 491, 56-65.
- 409 27. Baharian, S., Barakatt, M., Gignoux, C.R., Shringarpure, S., Errington, J., Blot, W.J.,
410 Bustamante, C.D., Kenny, E.E., Williams, S.M., Aldrich, M.C., et al. (2016). The Great
411 Migration and African-American Genomic Diversity. *PLoS Genet* 12, e1006059.
- 412 28. Mimno, D., Blei, D.M., and Engelhardt, B.E. (2015). Posterior predictive checks to quantify
413 lack-of-fit in admixture models of latent population structure. *Proceedings of the National*
414 *Academy of Sciences of the United States of America* 112, E3441-3450.
- 415 29. Bell, G.C., Caudle, K.E., Whirl-Carrillo, M., Gordon, R.J., Hikino, H., Prows, C.A., Gaedigk,
416 A., Agundez, J., Sadhasivam, S., Klein, T.E., et al. (2016). Clinical Pharmacogenetics
417 Implementation Consortium (CPIC) guideline for CYP2D6 genotype and use of
418 ondansetron and tropisetron. *Clin Pharmacol Ther*.
- 419 30. Hicks, J.K., Sangkuhl, K., Swen, J.J., Ellingrod, V.L., Muller, D.J., Shimoda, K., Bishop, J.R.,
420 Kharasch, E.D., Skaar, T.C., Gaedigk, A., et al. (2016). Clinical pharmacogenetics
421 implementation consortium guideline (CPIC) for CYP2D6 and CYP2C19 genotypes and
422 dosing of tricyclic antidepressants: 2016 update. *Clin Pharmacol Ther*.
- 423 31. Hicks, J.K., Bishop, J.R., Sangkuhl, K., Muller, D.J., Ji, Y., Leckband, S.G., Leeder, J.S.,
424 Graham, R.L., Chiulli, D.L., A, L.L., et al. (2015). Clinical Pharmacogenetics
425 Implementation Consortium (CPIC) Guideline for CYP2D6 and CYP2C19 Genotypes
426 and Dosing of Selective Serotonin Reuptake Inhibitors. *Clin Pharmacol Ther* 98, 127-
427 134.
- 428 32. Caudle, K.E., Rettie, A.E., Whirl-Carrillo, M., Smith, L.H., Mintzer, S., Lee, M.T., Klein, T.E.,
429 Callaghan, J.T., and Clinical Pharmacogenetics Implementation, C. (2014). Clinical
430 pharmacogenetics implementation consortium guidelines for CYP2C9 and HLA-B
431 genotypes and phenytoin dosing. *Clin Pharmacol Ther* 96, 542-548.
- 432 33. Johnson, J.A., Caudle, K.E., Gong, L., Whirl-Carrillo, M., Stein, C.M., Scott, S.A., Lee, M.T.,
433 Gage, B.F., Kimmel, S.E., Perera, M.A., et al. (2017). Clinical Pharmacogenetics
434 Implementation Consortium (CPIC) Guideline for Pharmacogenetics-Guided Warfarin
435 Dosing: 2017 Update. *Clin Pharmacol Ther*.
- 436 34. International Warfarin Pharmacogenetics, C., Klein, T.E., Altman, R.B., Eriksson, N., Gage,
437 B.F., Kimmel, S.E., Lee, M.T., Limdi, N.A., Page, D., Roden, D.M., et al. (2009).
438 Estimation of the warfarin dose with clinical and pharmacogenetic data. *The New*
439 *England journal of medicine* 360, 753-764.
- 440 35. Foster, M.W., Sharp, R.R., and Mulvihill, J.J. (2001). Pharmacogenetics, Race, and
441 Ethnicity: Social Identities and Individualized Medical Care. *Therapeutic Drug Monitoring*
442 23, 232-238.

- 443 36. Yen-Revollo, J.L., Auman, J.T., and McLeod, H.L. (2008). Race does not explain genetic
444 heterogeneity in pharmacogenomic pathways. *Pharmacogenomics* 9, 1639-1645.
- 445 37. Braun, L., Fausto-Sterling, A., Fullwiley, D., Hammonds, E.M., Nelson, A., Quivers, W.,
446 Reverby, S.M., and Shields, A.E. (2007). Racial categories in medical practice: how
447 useful are they? *PLoS medicine* 4, e271.
- 448 38. Ortega, V.E., and Meyers, D.A. (2014). Pharmacogenetics: implications of race and ethnicity
449 on defining genetic profiles for personalized medicine. *J Allergy Clin Immunol* 133, 16-
450 26.
- 451 39. Bamshad, M., Wooding, S., Salisbury, B.A., and Stephens, J.C. (2004). Deconstructing the
452 relationship between genetics and race. *Nat Rev Genet* 5, 598-609.
- 453 40. Urban, T.J. (2010). Race, ethnicity, ancestry, and pharmacogenetics. *Mt Sinai J Med* 77,
454 133-139.
- 455 41. Risch, N., Burchard, E., Ziv, E., and Tang, H. (2002). Categorization of humans in
456 biomedical research: genes, race and disease. *Genome Biology* 3.
- 457 42. Bamshad, M.J., Wooding, S., Watkins, W.S., Ostler, C.T., Batzer, M.A., and Jorde, L.B.
458 (2003). Human population genetic structure and inference of group membership. *Am J*
459 *Hum Genet* 72, 578-589.
- 460 43. Bustamante, C.D., Burchard, E.G., and De la Vega, F.M. (2011). Genomics for the world.
461 *Nature* 475, 163-165.

463

464

465 **Figures:**



466

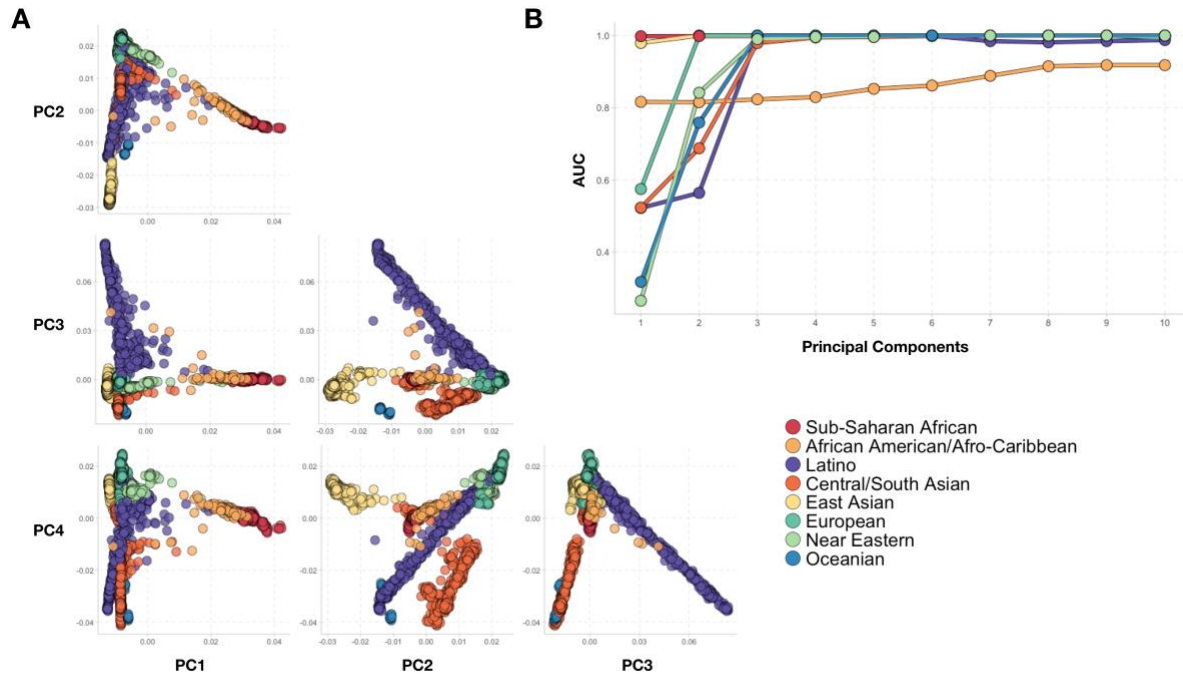
467

468 **Figure 1: Map of geographical boundaries included in each geographical population group.**

469 Group boundaries for the seven geographical groups fall predominantly along national
470 boundaries to aid the assignment of group membership. The two admixed groups of African
471 American/Afro-Caribbean and Latino are not shown on this figure as the map indicates the
472 borders of each geographical group based on the location of genetic ancestors pre-Diaspora and
473 pre-colonization, which cannot be applied to the two admixed groups. It should also be
474 recognized that, due to the large geographical areas covered by each group, a single group does
475 not accurately represent the large amount of genetic diversity found in that one region.

476

477



478

479 **Figure 2: Principal component analysis comparing genetic distances of populations with**

480 **close geographic proximity using 1000 Genomes participants. (A) The genetic gradient**

481 between populations is illustrated along PCs 1 vs 2 and PCs 3 vs 4, showing that, while

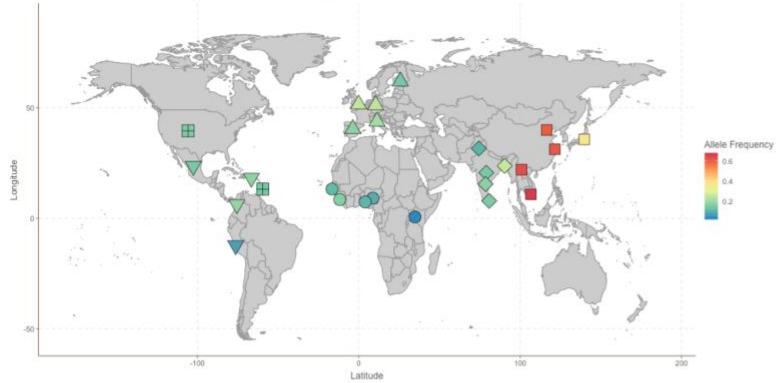
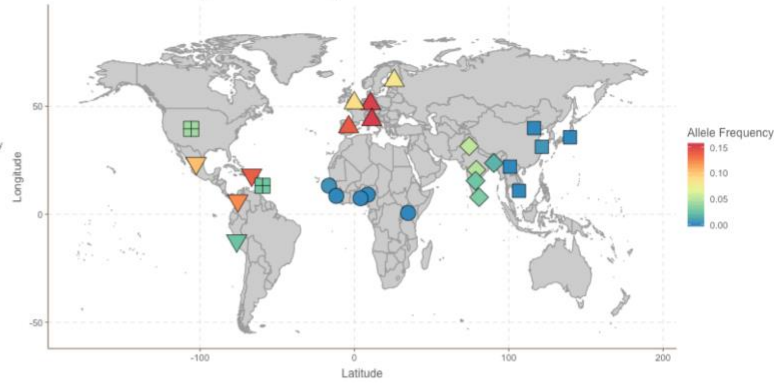
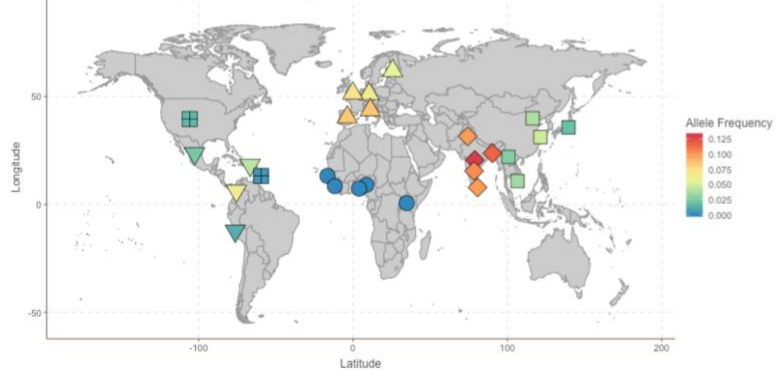
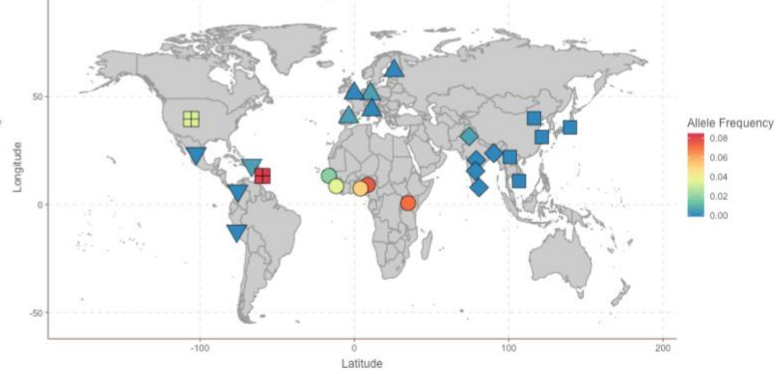
482 completely discrete population boundaries are challenging, the groupings proposed here provide

483 a statistically robust grouping. (B) AUCs of logistic regression to predict cluster membership,

484 showing high degree of population structure. Note that, because none of the 1000 Genomes

485 populations fall into the American (AME) group, no reference data were available to include this

486 group in the analysis.

A. *CYP2D6*10* (rs1065852)B. *CYP2C9*2* (rs1799853)C. *CYP2C9*3* (rs1057910)D. *CYP2C9*8* (rs7900194)

487

488 **Figure 3: Maps illustrating how the proposed biogeographical grouping system can be used to illustrate the variability in**
 489 **global frequencies of key pharmacogenetic alleles.** Allele frequencies from 1000 Genomes are shown across global populations for
 490 (A) *CYP2D6*10*, (B) *CYP2C9*2*, (C) *CYP2C9*3* and (D) *CYP2C9*8*.

